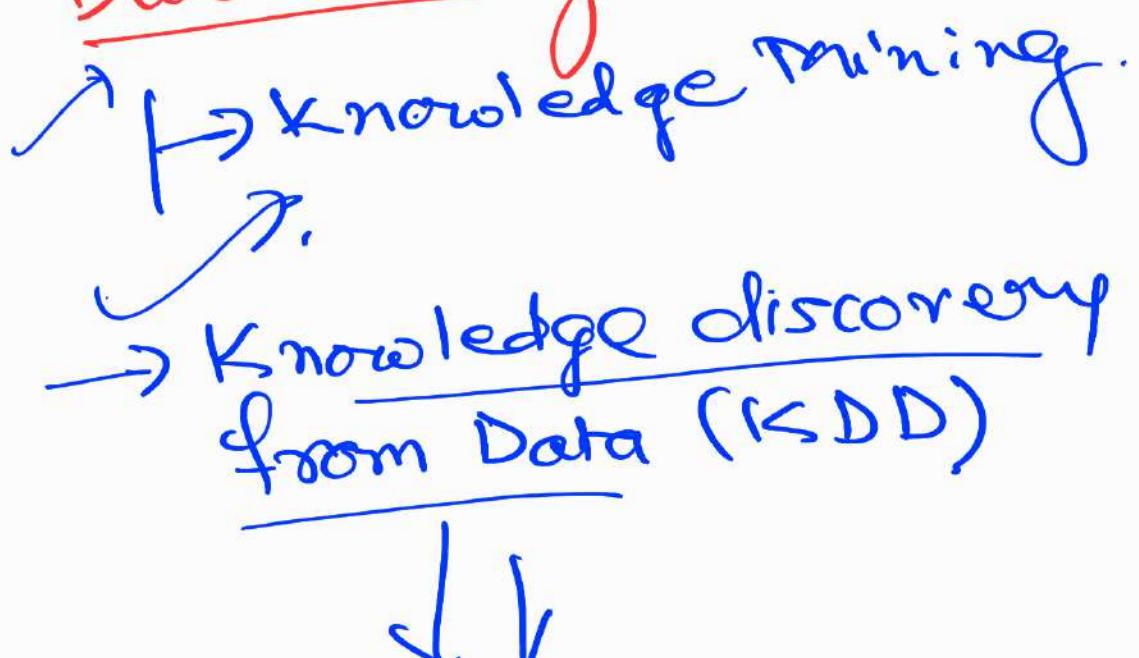


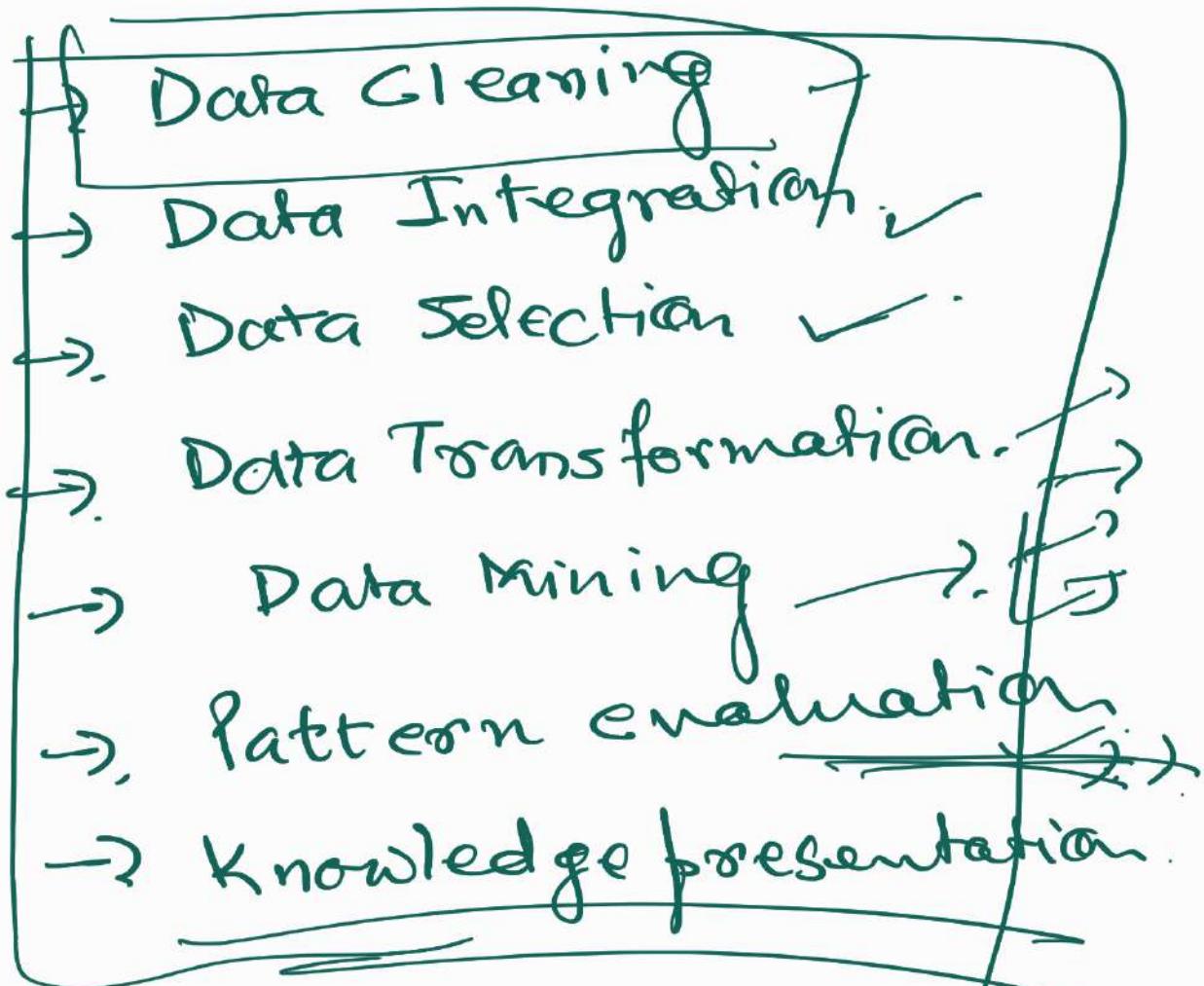
m

Data Mining & Warehousing

Data Mining:-



Extracting the knowledge
from the large amount
of Data.



- Temporal Databases
- Sequence Database
- Time-series Database

Temporal Databases

→ Time related attributes.

→ timestamp

Sequence Database :-

| → sequence of events or activity happens.

Time Series Database :-

| → sequence of values or events over

repeated measurement of time (minute / second / millisecond) daily / weekly . etc

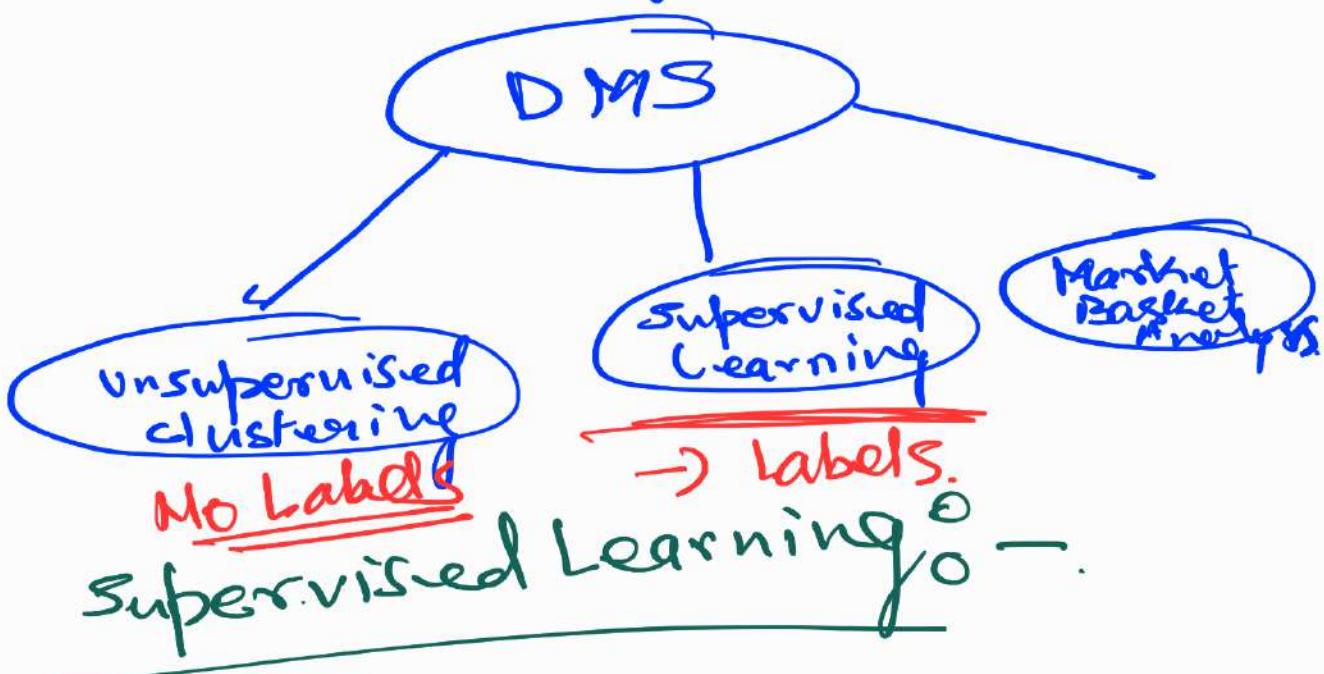
Ex:-

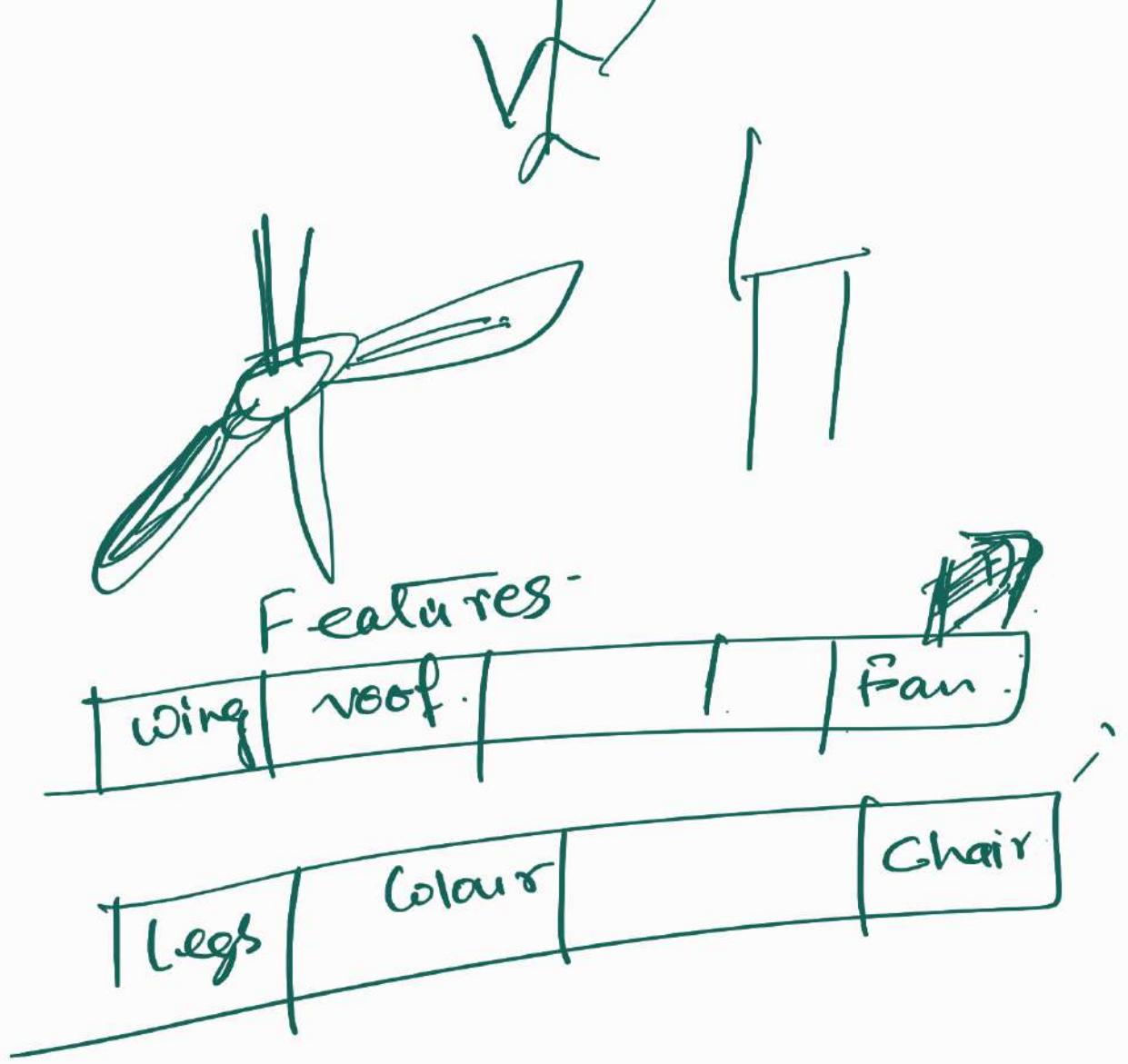
Data Mining

↓
→ Data mining (KDD)
is the process of analyzing
data from different ^{perspective}
summarizing it into
useful information that
can be used to
→ increase revenue.
→ cut ~~out~~ cost

or
Both -

Data Mining strategies

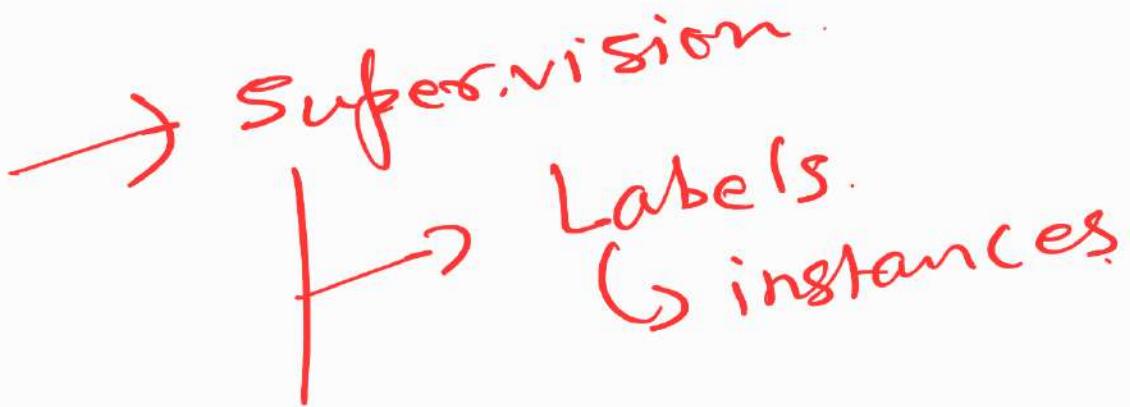




- When we are young
- we use induction to form basic concept.
- animals
plants
buildings
- mind choose what

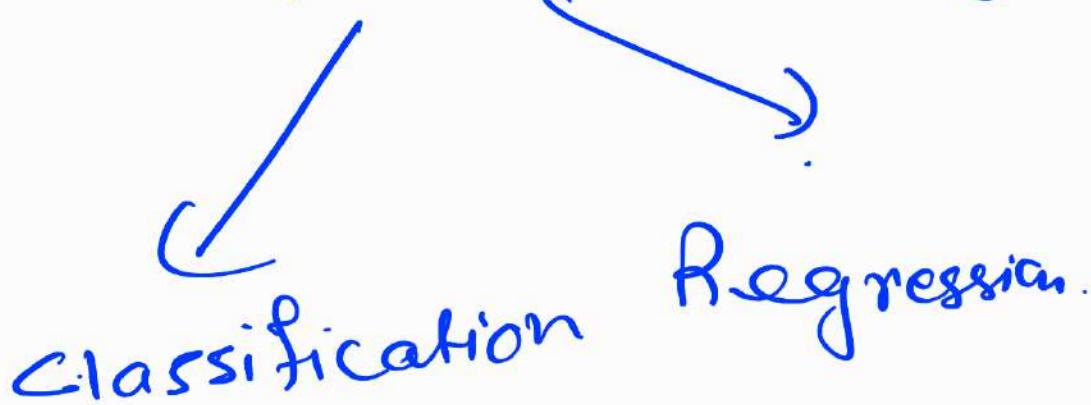
it believe to be the defining concept feature

↳ form our own classification Model.



→ induction based learning.

Supervised Learning



Classification -

→ the output is

→ a certain defined

having certain defined
labels (discrete values)

→ Binary or }
Multiclass }

Op → Yes | No

Regression :- → Red, Green
Blue.

→ here the output is
having continuous values

Ex :-			
Temp	Press	Humidity	Wind speed
10.34	924.07	—	?

→ Predict the sale.
→ Price of house
,, , " share

Unsupervised Clustering :-

→ it builds model
from data without

Predefined classes

→ Data instances are grouped together based on the similarity scheme defined by clustering system.

→ Clustering

→ grouping the data into clusters, so that objects within a cluster have high similarity ~~or~~

→ but are very dissimilar to objects in other clusters.



- Dissimilarity is assessed based on the attribute values that are describing the object
- Distance measure is generally used.



First we apply clustering algo

↳ Labelled it.

↓
Labelled Data.

→ Data Mining
→ steps.

→ = types.
of DM Alg.

→ Market Basket
Analysis

Mobile → Mobile Cover
Mobile → Tempered Glass

→ Rules
→ Sales
→ Association Rule

A → B
item item

Mobile → Tempered Glass.

↓

Arrange.

Mobile, Tempered Glass

10:30 — hrs.

Data Mining

Mean :-

Take few numbers

13, 18, 13, 14, 13, 16, 14, 21, 13

mean = ?

$$\frac{\text{Sum of num}}{\text{Total No.}} = \frac{135}{9}$$

Median :-

13, 13, 13, 13, 14, 14, 16, 18, 21

find the median ?

Ans \Rightarrow 14 8

Median \Rightarrow it is the middle value of the data set when arranged in

it has been arranged in
order.

Ex 2 :-

12, 18, 16, 21, 10, 13, 17, 19

if the data set is even.

⇒ median is the average
of two middle values

→ 10, 12, 13, 16, 17, 18, 19, 21

$$\Rightarrow \frac{4^{\text{th}} \text{ data val} + 5^{\text{th}} \text{ data val}}{2}$$

$$\Rightarrow \frac{16+17}{2} = 16.5$$

Mode :-

48, 44, 48, 45, 42, 49, 48

⇒ data values in the
dataset that occurs

most often.

Ex 8 -

dataset is
given like this

age	frequency
1-5	200
5-15	450
15-20	300
20-50	1500
50-80	700
80-110	44

Calculate the Approx median value of the data.

$$\text{median} = L_1 + \left(\frac{\frac{N}{2} - (\sum \text{freq})_1}{\text{freq}_{\text{median}}} \right) \text{width}$$

boundaries of

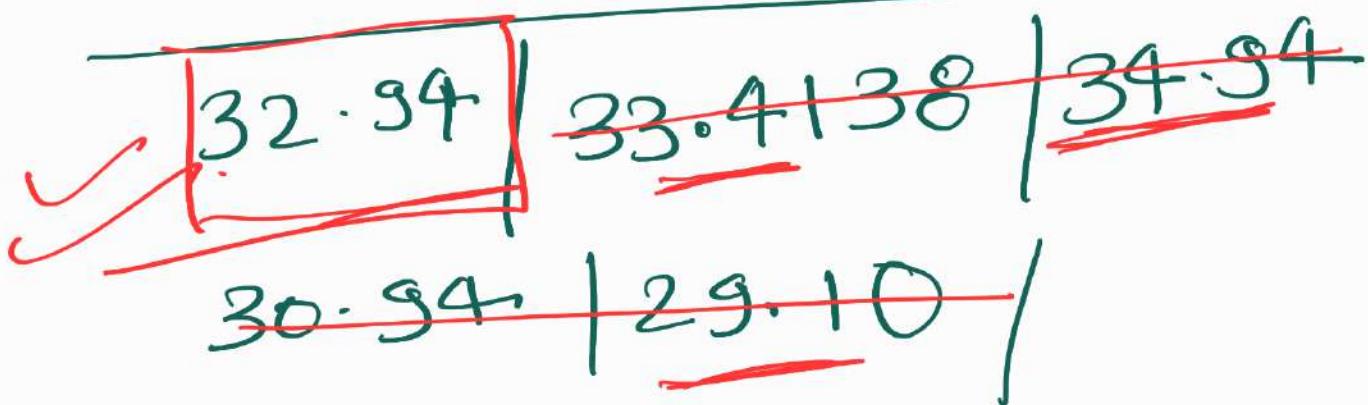
$L_1 \rightarrow$ Lower boundary of the median interval.

$N \rightarrow$ is the no. of values in the data set

$(\sum \text{freq})_l \rightarrow$ sum of frequencies of all the intervals that are lower than the median interval

$\text{freq}_{\text{median}} \rightarrow$ is the freq of the median interval

width \rightarrow it is the width of the median interval



$$N = \underline{3194} \quad \text{median} \Rightarrow 1597$$

$$L_1 = 20$$

$$\sum \text{freq}_i = 150$$

$$\text{freq}_{\text{median}} = 150$$

$$\text{width} = 30^\circ$$

$$\text{median} = \underline{\underline{32.94}}$$

Data Preprocessing

→ Data Cleaning

- incomplete,
noisy, inconsistent
- fill in the missing values
- ~~Smooth out the noise~~
while identifying
outliers
- Correct Inconsistency
in the data

Missing Values

Ignore the tuple

Age	Address	Gender	F3	F4	F5	Glass label
f ₁	f ₂	f ₃	f ₄	f _n		
X	C	C	C	X	C	X.

(A)

X	X	X	X	✓	✓	X	X	✓
---	---	---	---	---	---	---	---	---

(B)

↙ Not present

→ Fill in the missing values
manually :-

→ Use a global constant to
fill in the missing
value

Ex unknown, 0, or
etc

~~➤~~ use the attribute mean to fill in the missing value.

f_i	Income	Glass
34	100	
26	200	
33	mean	
28	300	
33	150	
11	200	

~~➤~~ Use the attribute mean for all samples belonging to the same class as the given tuple.

f_i	Income	f_n	Glass
	100	A	A

34	200		A
26	300		A
31	?	meanA	A
29	100	↑	?
38	50	↓	B
11	meanB		B
12	20		B
19			B

⇒ Use the most probable values to fill in the missing value.

→ Ex - Regression
 Bayesian
 Decision Tree

Standard Deviation

1, 2, 3, 4, 5, 6

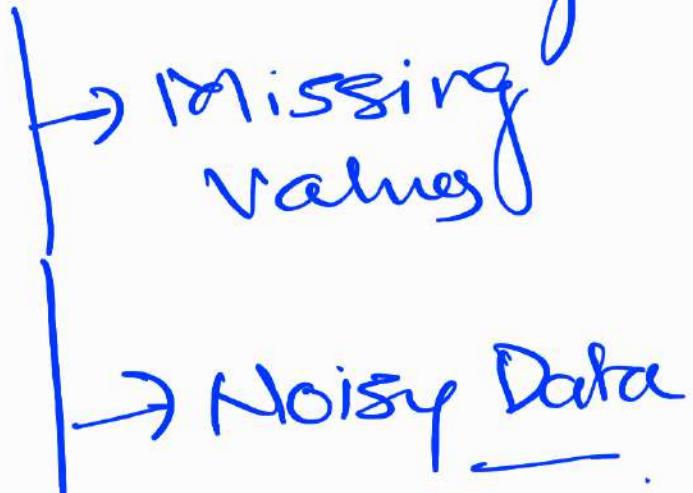
S. D. $\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}}$

mean = 3.5

SD = 1.7



→ Data Cleaning



Noisy Data :-

↳ random error

- Binning
- Regression
- Clustering

Binning :-

It smooths data values by consulting its neighbourhood. i.e. the value around it.

Ex :- 4, 8, 15, [21, 21, 24], 25, 28, 34

→ Put the data values into Buckets / Bins.

→ Local smoothing

↳ Smoothing by bin

→ smoothing by means

→ smoothing by bin boundaries

Bin 1: 4, 8, 15

Bin 2: 21, 21, 24

Bin 3: 25, 28, 34

smoothing by bin means

Bin 1: 9, 9, 9

Bin 2: 22, 23, 22

Bin 3: 29, 29, 29

smoothing by bin boundaries.

Bin 1: 4, 4, 15

min = 4 marks

Bin 2: 21, 21, 24

min = 21 marks

Bin 3: 25, 25, 34

max = 34

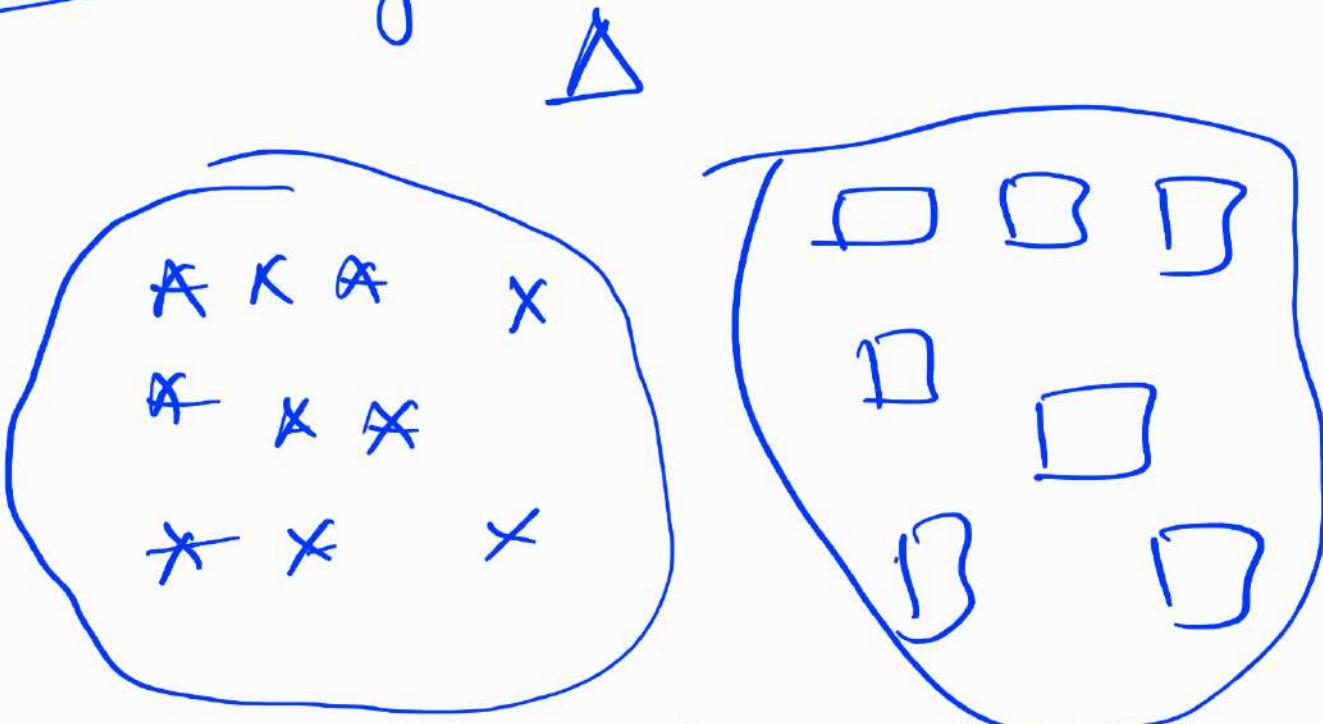
min & Max values are

bin boundaries identified. \Rightarrow bin boundaries
 \rightarrow Each bin values is then replaced by the closest boundary value

Regression -

\rightarrow data can be smoothed by fitting it into a function

Clustering -



G1

G2

Outliers

→ They fall outside the set of clusters.

Data Cleaning :-

→ discrepancy detection

→ metadata

→ Data scrubbing tools

↓
use domain knowledge
to detect & correct
data.

⇒ Data Auditing tool.



it finds rules &
relationship



Data Mining Techniques

Correlation Analysis -

↓
↓
relationship b/w
Numerical
Attributes



Pearson Correlation.

$$\rho_{A,B} = \frac{\sum_{i=1}^N (a_i - \bar{A})(b_i - \bar{B})}{\sqrt{\sum_{i=1}^N (a_i - \bar{A})^2 \sum_{i=1}^N (b_i - \bar{B})^2}}$$

$\nabla \rightarrow$ No. of tuples

a_i^j & b_i^j are respective values
of the tuple.

$$-1 \leq r_{A,B} \leq +1 .$$

-ve value \Rightarrow Negatively correlated

+ve value \Rightarrow Positively correlated

0 \Rightarrow A & B are independent & no correlation

Tree Height	Trunk Diameter
4	r_{x}

35	8
49	9
27	7
33	6
60	13
21	7
45	11
51	12

Find the correlation b/w
Tree Height & Trunk
Diameter.

Ans 8 - $\simeq 0.89$

+ positively
correlated

Data Transformation

- ④ \Rightarrow transformed or consolidated into

One forms appropriate for mining.

\doteq → smoothing.
→ removed of
noise
→ Binning
→ Regression
→ Clustering

Σ → Aggregation -
| → Aggregation opp.
| → daily → weekly

3 → Generalization

→ We will pass
the low level
data by the higher
level concept through
the concept hierarchy

Ex: - city can →
state

↳ → Normalization ~

↳ Data are scaled
so as to fall with
in a specified range

→ 0 to +1

OR

0 to 1.

↳ → Attribute Construction

↳ new attributes
are created from

are constructed from
the given set of
attributes.



Data Transformation

- Smoothing
- Aggregation
- Generalization
- Normalization
- Attribute construction

Normalization Techniques

- Min-max normalization

$\underline{\min_A} \rightarrow$ min value of an attribute A.

$\underline{\max_A} \rightarrow$ max. value of an attribute A.



In the new range.

$\underline{\text{new_min}_A} \rightarrow$ New min val of an attribute A.

$\text{new_max}_A \rightarrow$ new max value
of an
attribute A.

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new_max}_A - \text{new_min}_A) \rightarrow + \text{new_min}_A$$

Ex 2 - min & max value for
attribute income is
₹ 12000/- & ₹ 98000/-
map into range.
[0.0 to 1.0]

to find the value of
₹ 73600/-

$$0.71 / 0.716$$

Score normalization

Zero-mean normalization

value of an attribute A
are normalized



based on mean &
standard deviation of A.

initial value $\rightarrow v$.

After normalization $\rightarrow v'$

$$v' = \frac{v - \bar{A}}{\sigma_A}$$

$\bar{A} \rightarrow$ mean

$\sigma_A \rightarrow$ standard deviation

Ex

Attribute income

mean \rightarrow 75.5

SD \rightarrow $\sqrt{160}$ or 12.7

with z-score non " find
the value of $\mathbb{P}(Z \geq 1.225)$

Ans 1.225

Boxplot :-

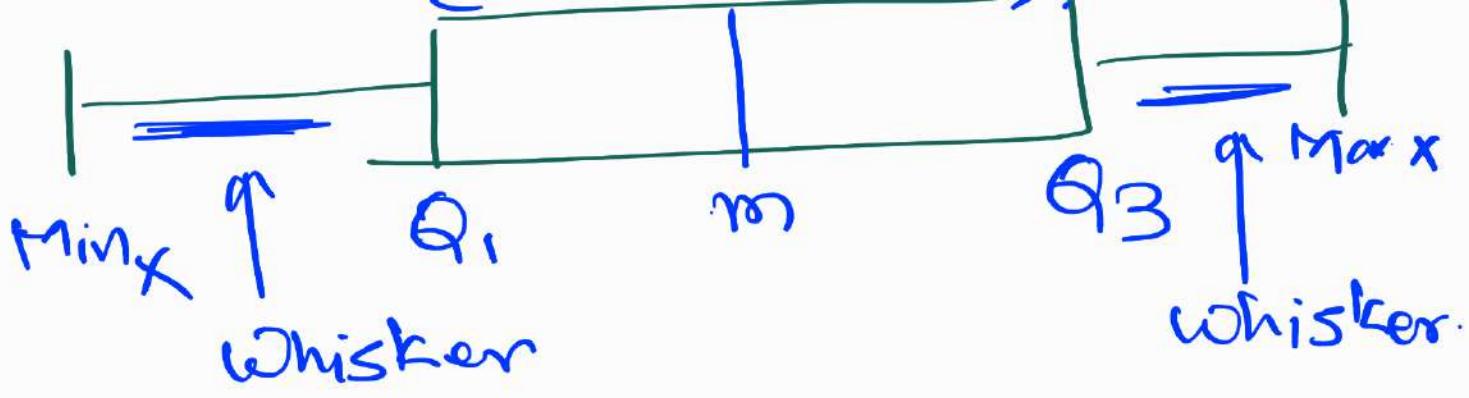
76, 79, 76, 74, 75, 71, 85, 82,
82, 79, 81

\rightarrow Arrange the values in
the ascending order.

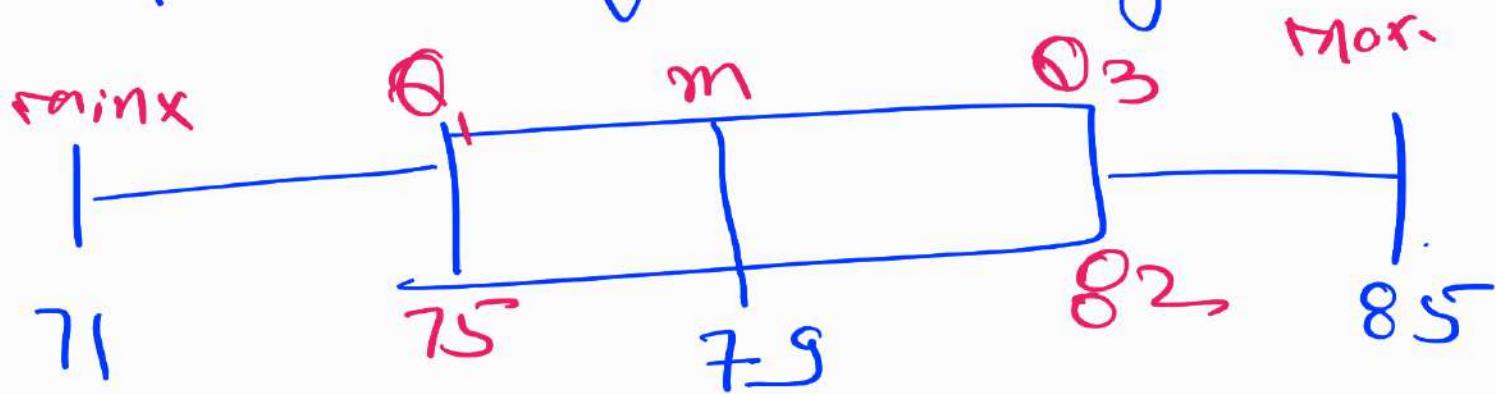
71, 74, 75, 76, 76, 79, 79, 81, 82, 82,
85

\rightarrow It gives 5 numbers
summary of the data
distribution

IQR



IQR: Interquartile Range.



71, 74, 75, 76, 76

79, 81, 82, 82, 85

α

Pearson correlations -

$$-1 \leq r_{AB} \leq 1$$

$\rightarrow 1 \Rightarrow$

$\times 1 \Rightarrow$

$0 \Rightarrow$

χ^2 (Chi square Test)

→ It finds the correlation between two categorical attributes

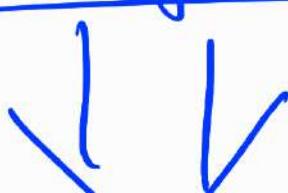
lets say those attributes are A & B.

→ A \Rightarrow 'c' distinct values
 $a_1, a_2, a_3, \dots, a_c$

→ B \Rightarrow 'd' distinct values
 $b_1, b_2, b_3, \dots, b_d$

represented in the form of

'Contingency Table'



'c' values of 'A' making up the col.

'r' value of 'B' making up the rows.

$(A_i, B_j) \rightarrow$ denote an event that

$$A \rightarrow a_i \\ B \Rightarrow b_j$$

$$A = a_i \\ B = b_j$$

→ In contingency table we represent different joint event in diff cells

$$\chi^2 \text{ value} = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

where O_{ij} is the observed frequency (actual count) of the joint-event (A_i, B_j)

E_{ij} is the expected frequency of (A_i, B_j)

$$E_{ij} = \frac{\text{Count}(A=a_i) \times \text{Count}(B=b_j)}{N}$$

$N \rightarrow$ No. of data tuples

\Leftrightarrow No. of tuples having value a_i for A.

$\rightarrow b_j$ for B

Σx_i^o

observed frequency

	male	female	Total
--	------	--------	-------

		<u>250</u>	(90)	<u>200</u>	(360)	450
		<u>50</u>	(210)	<u>1000</u>	(840)	1050
		<u>300.</u>		<u>1200.</u>		1500
						Expected freq

$$e_{11} = \frac{\text{count(male)} \times \text{count(fiction)}}{N}$$

(male, fiction)

$$\approx \frac{300 \times 450}{1500} = 90$$

e
(fiction, female)

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} +$$

~~$$\frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840}$$~~

$$= 507$$

Q. 508

$H_0 \rightarrow$ There is no relationship
between categorical variables
A & categorical variable
B.

$H_1 \rightarrow$ There is some relation
ship between these two
categorical variables.

Step ②

Find the degree of
freedom

$$(R-1) \times (C-1)$$

$$\text{dof} = (2-1) \times (2-1) \\ = 1$$

Step ③

significance level.

1% or 5%

χ^2 , dof, SL
value, —, —

→ If the ~~value~~ value in the chart is less than ~~the~~ calculated χ^2 then we accept H_1

→ If the calculated χ^2 is less than we accept H_0

χ^2

Data Warehouse

- OLTP
- OLAP Systems.

OLTP :- we are using in day to day operation.

OLAP :- we put historical data after integration/ transformation and various other operations.

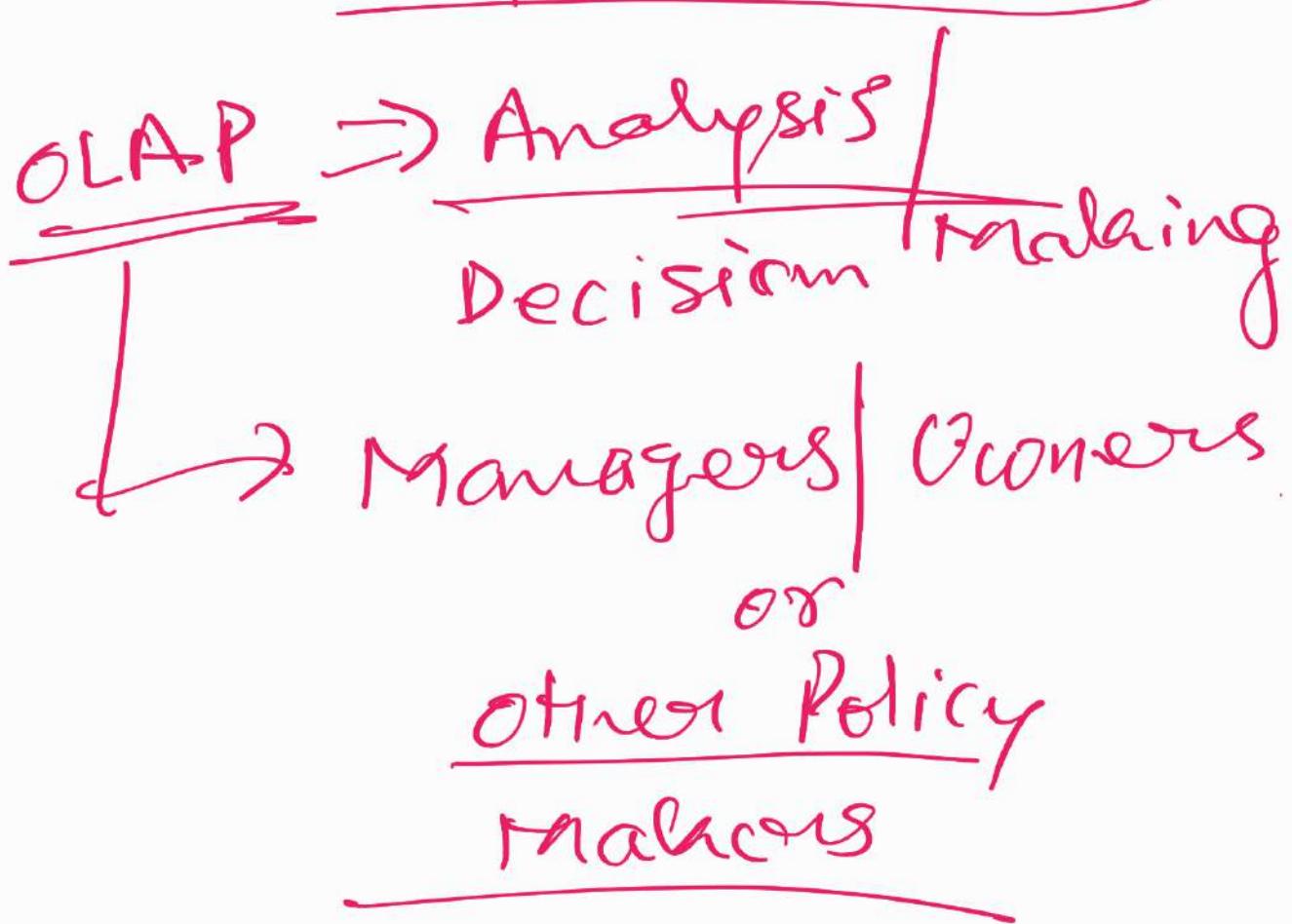
→ On-line analytical processing

→ 50 Transactions
G. ≈ 1 sec

$$1 Tx \rightarrow \frac{1}{50}$$

$$1 Gross \rightarrow \frac{1}{50} \times \frac{100000000}{\cancel{1000000}}$$

→
How much hours?
→ 55-5 hrs.



Total Revenue Generated
in Partner by Airtel or
Other



Data warehouse :-

- to systematically organize
- understand
- Use their data to make strategic decisions

William H. Inman :-

- Subject Oriented
 - Integrated
 - Time-variant
 - Non-volatile.
- in support of

management decision making process

Subject Oriented :-

↳ Centered around
major subjects, ex:-
Sales, Suppliers

Integrated :-

→ Integrating data
from multiple sources.
These same sources may
be homogeneous or
heterogeneous.

Time-Variant

→ data is stored from
the historical perspective

→ Every key structure in DW
contain either implicitly

or explicitly an element of time.

Nonvolatile -

→ we perform two operations

 → initial Loading

 → Access of data

Major Difference b/w OLTP & OLAP

OLTP → day to day operat.

OLAP → used to key decision making

Data Content -

OLTP → Current data,
 |
 | detailed

OLAP → large amount of historical, aggregated

data

Database Design

OLTP → ER models

OLAP → star or snowflake.

fact-constellation

Access Pattern

OLTP → short, atomic transactions

OLAP → read only operations

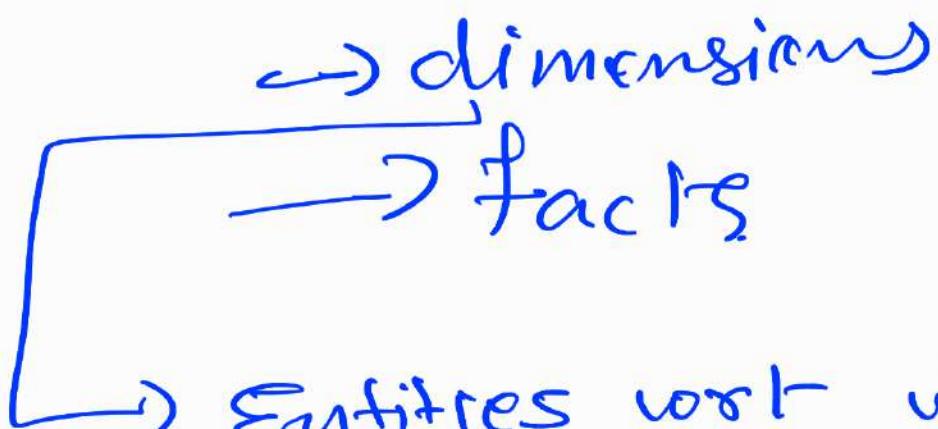
Data Cubes -

→ Multidimensional Data models.



Data Cubes

A data cube allows data to be modeled and viewed in multiple dimensions.



→ Entities over which an organization wants to keep records

. For Sales subject

dim → time, items, branch, location etc.

→ Each dim " has table associated with it



dimension table.

→ The subject around

which multi-dimⁿ
data model is organized
is represented by
FACT table.

Ex 8 -
2 D - view of sales
Data
dim → time & item

Sales from branches in
Patna

location = "Patna"

item(type)

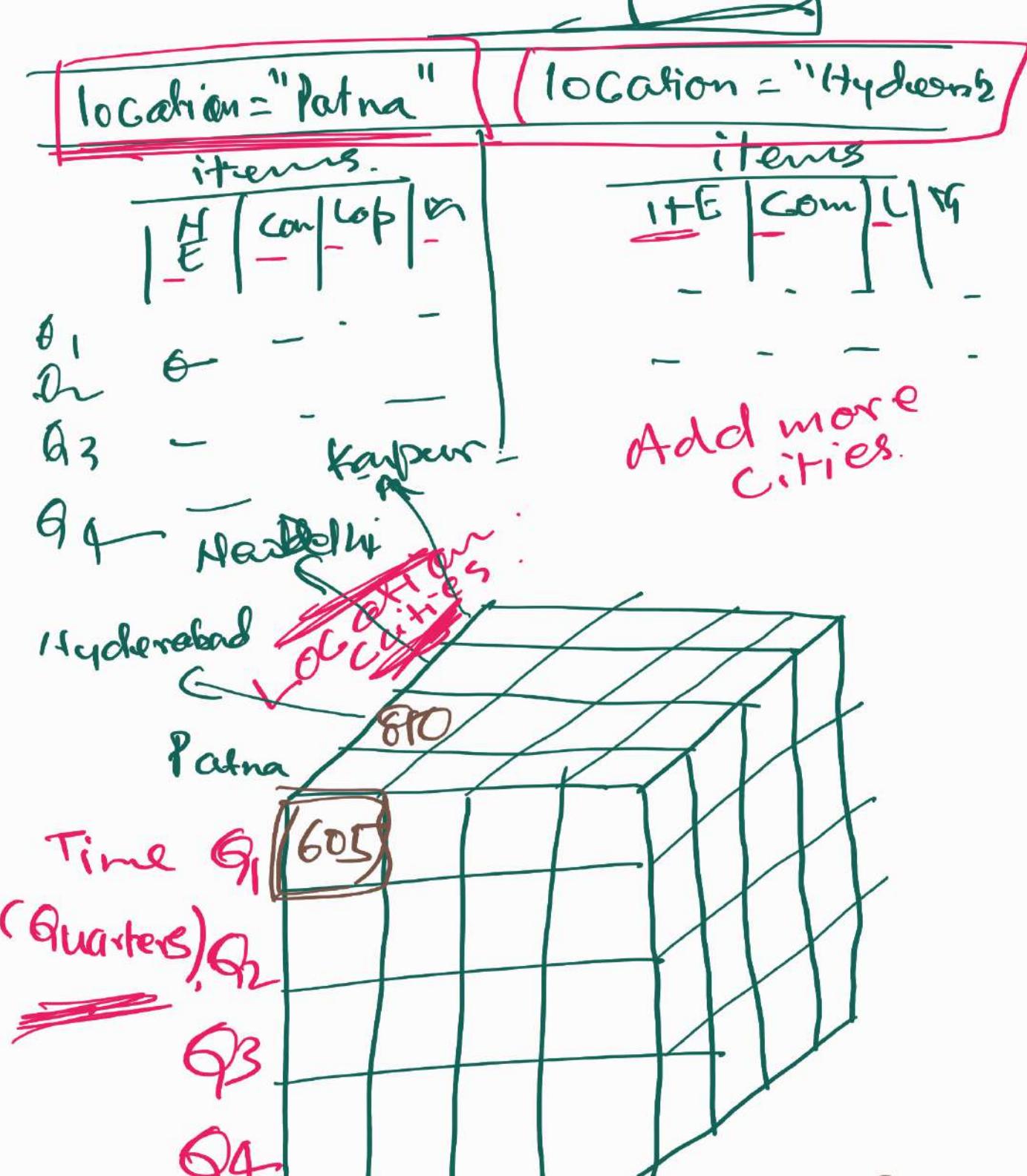
time(Quarter)	Home Enter		computer	Laptop	Mobile
	Q1	Q2	609	961	106
			514	102	-
			623	103	107
					-

Q3

Q4 691

104

→ Can organize the data
in 3D



Phone
Enter

Laptop

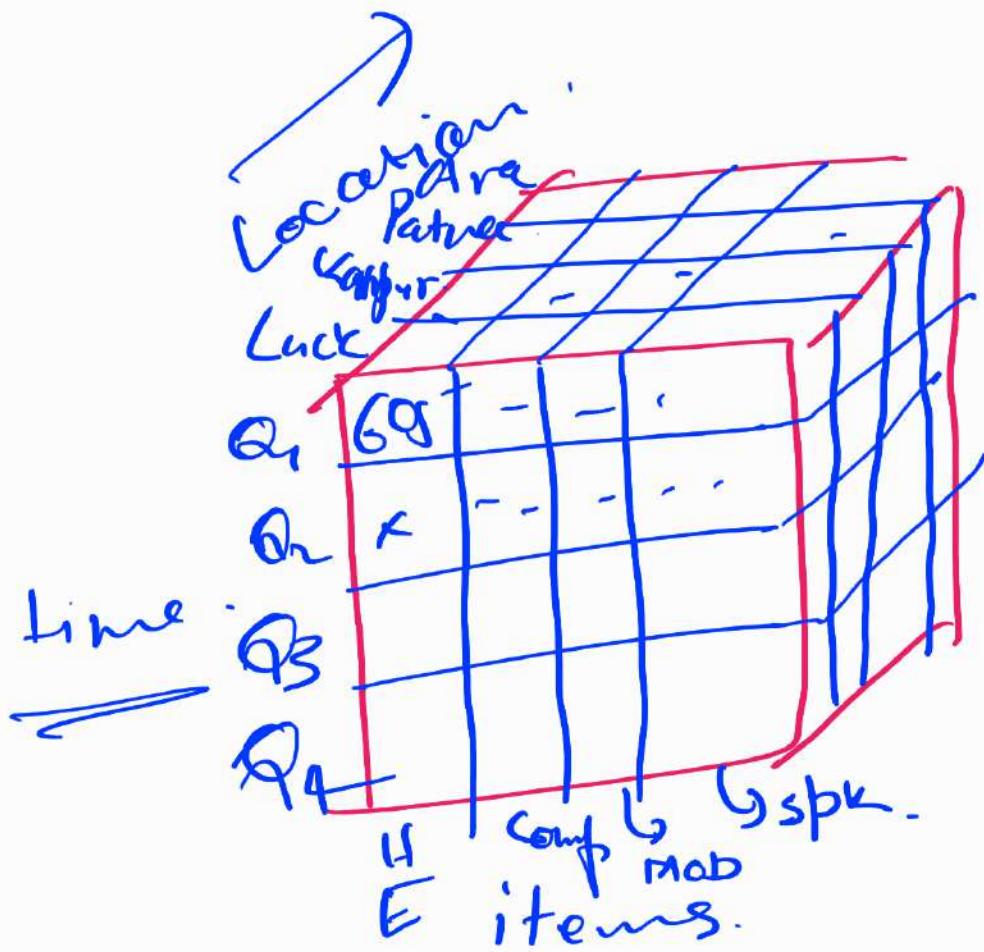
Mobile.

Comp

Items (types).



Data Cube :-



④ Four Dimⁿ

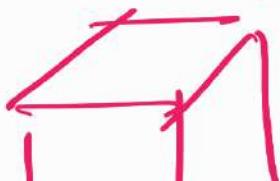
→ As a series of 3D cubes

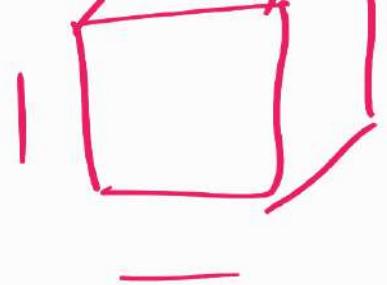
qth dim Supplier.
time, item, location .

Supplier = sup1

Supplier = sup2

Supplier = sup3





* If we have n -D data then we can represent it as a series of $(n-1)$ -D cubes

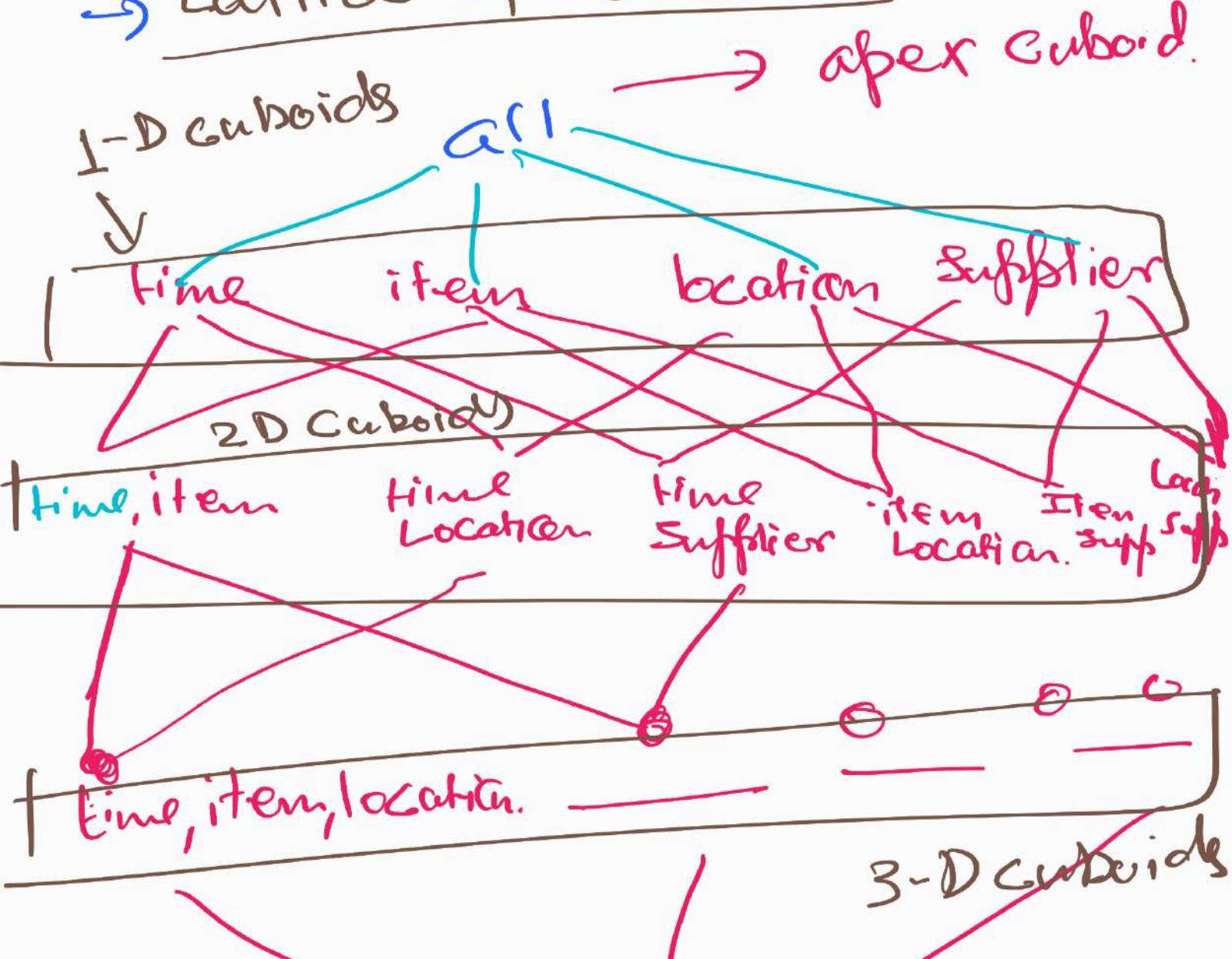


→ Lattice of cuboids :-

1-D cuboids

→ apex cuboid.

$a^{(1)}$



time, item, location.
Base Cuboid. Supplier.

Lattice of cuboids

Concept Hierarchies

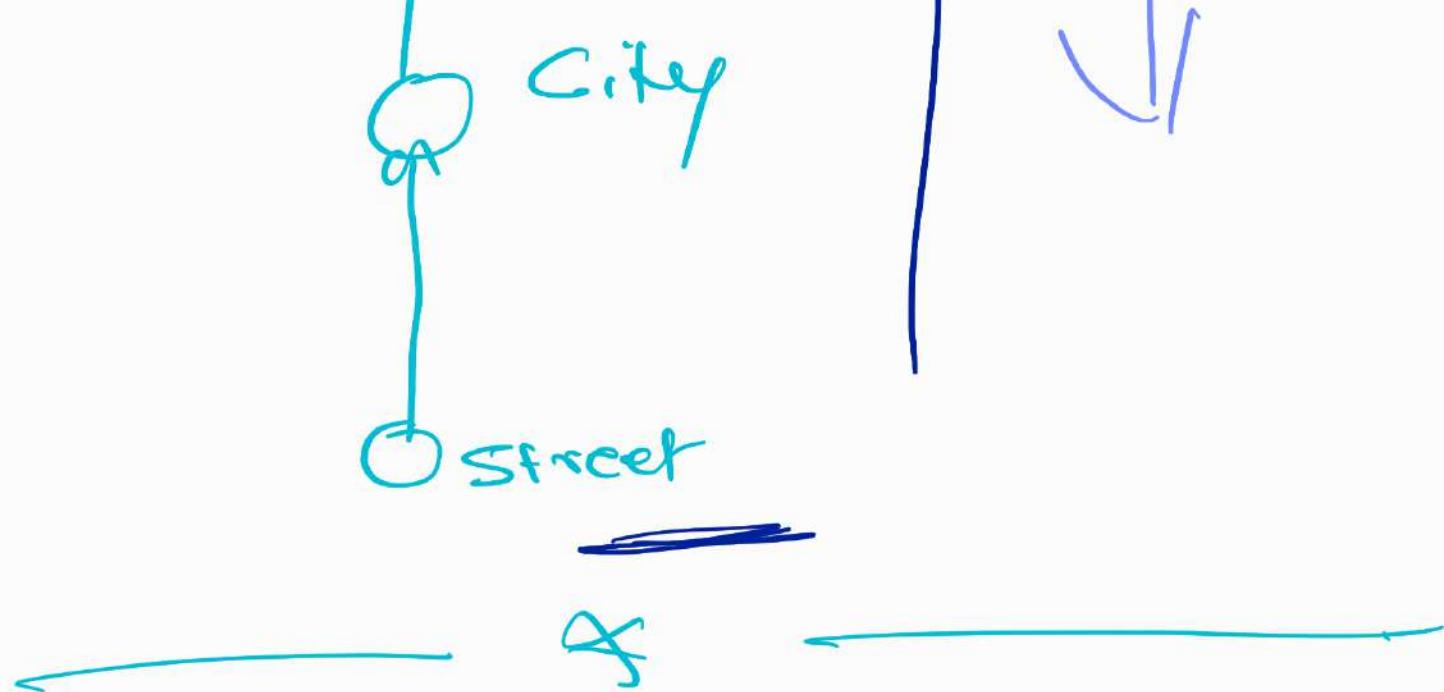
→ It defines a sequence of mapping from a set of low level concepts to higher level or more general concepts

Country



state





Schemas for multidimensional Databases

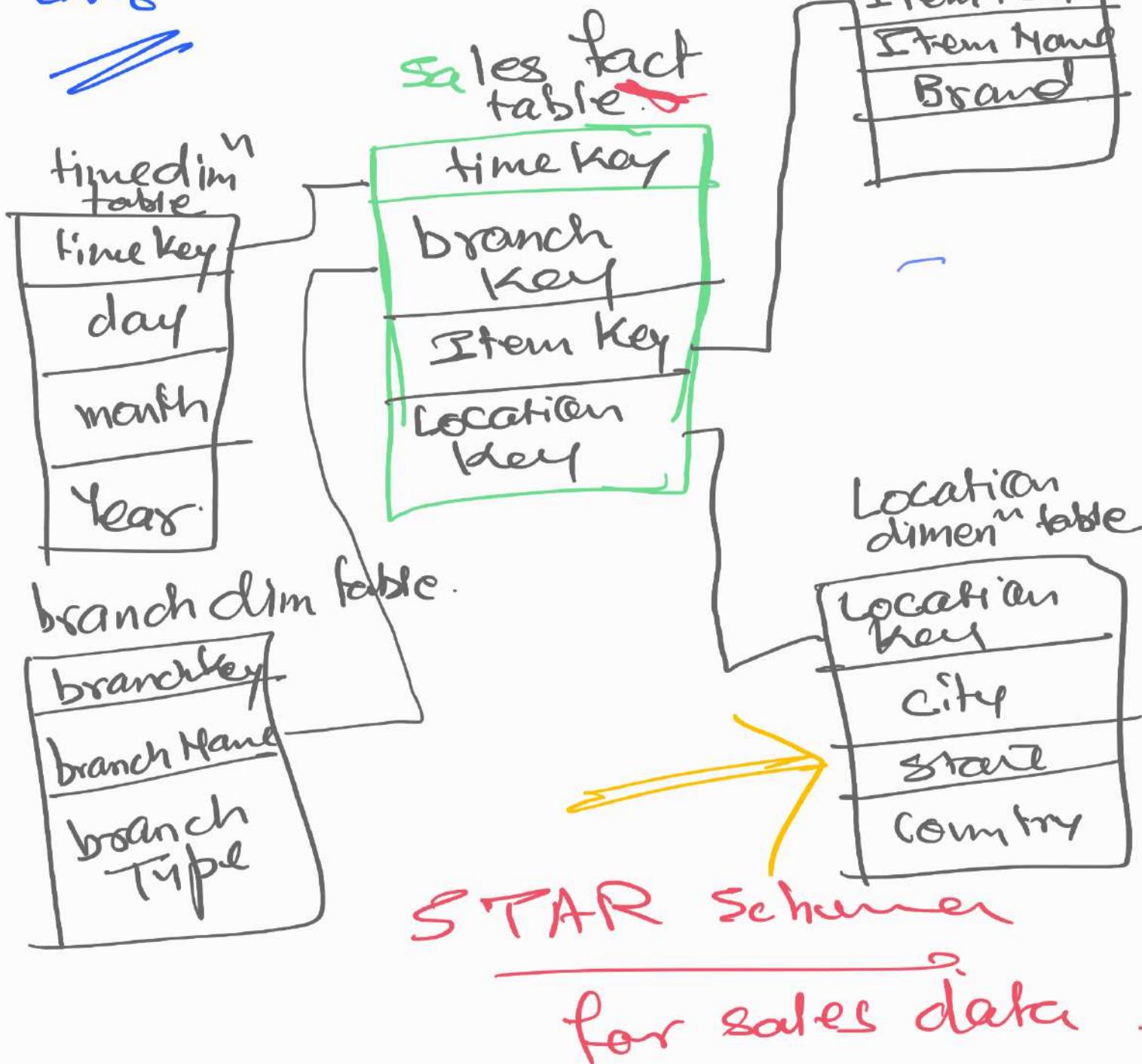
- Star schema
- Snowflake schema
- Fact constellation schema.

Star schema :-

- It has large central table (fact table) ~~and~~ containing the bulk of data with no redundancy
- Set of smaller attendant tables (dimension tables)

→ one for each dimension

~~EXO~~



Snowflake schema's -

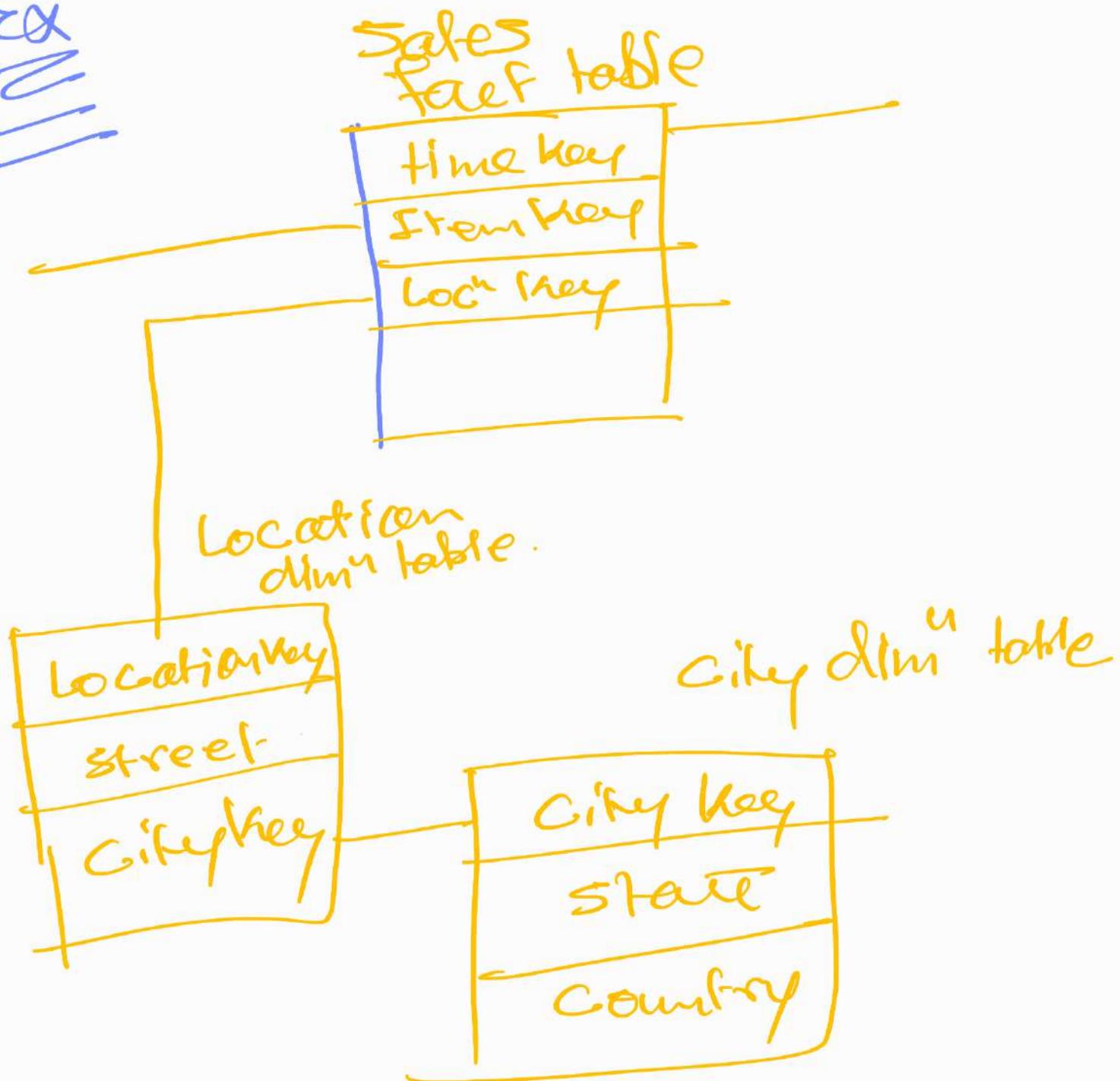
* Variant of STAR .

** In a dimension table one

some
normalized.

→ Further splitting of
data.

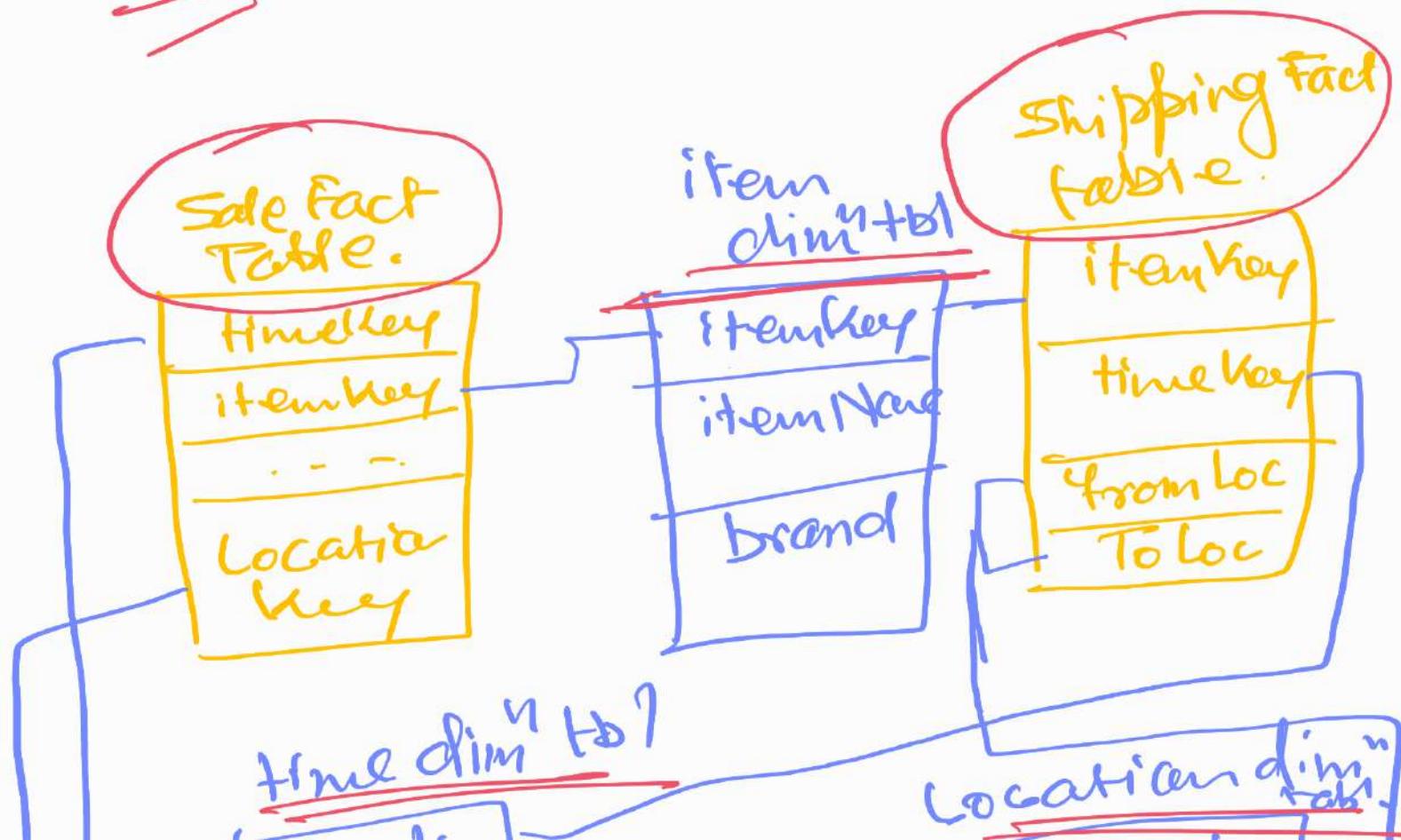
Ex
|||

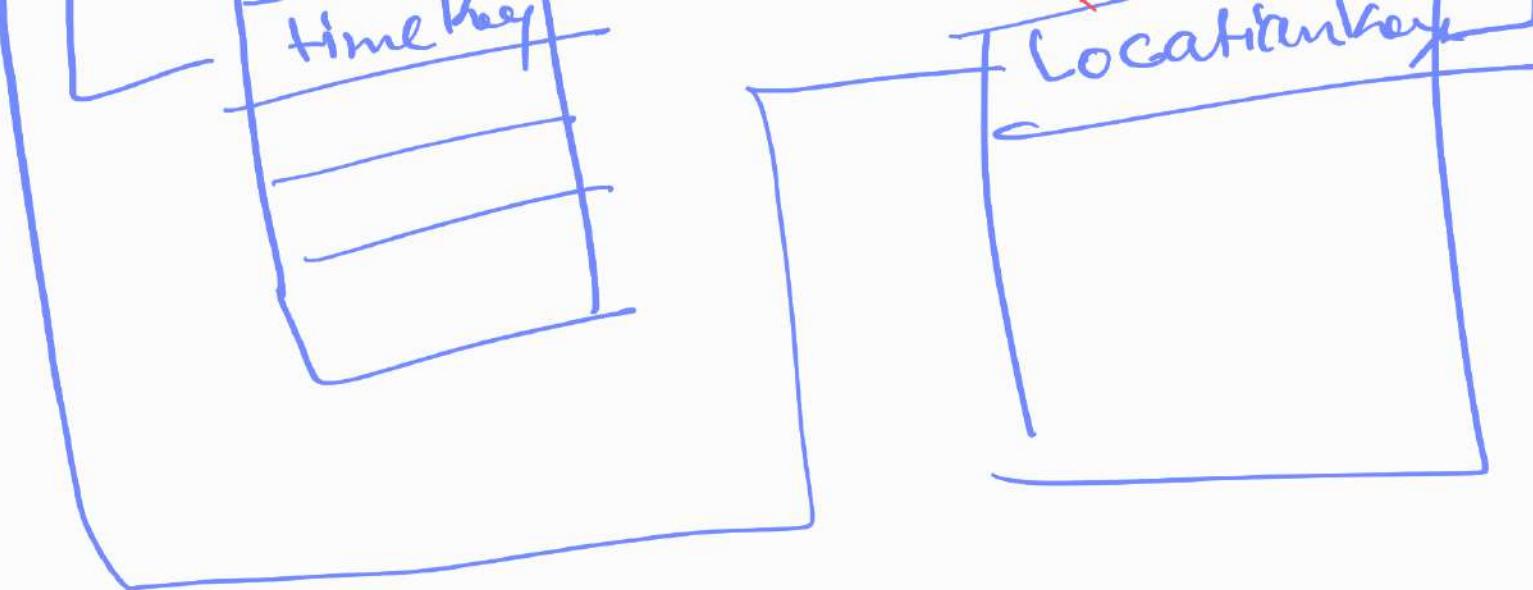


Fact Constellation

- Galaxy schema
- sophisticated | complex off
- Has multiple fact tables that
 - ↓
 - share
 - ↓
 - dimension tables

Ex 8



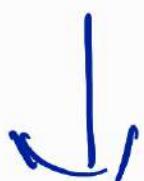


Fact constellation
schema for sales & shipping

OLAP operations :-

→ Roll-up :-

→ performs aggregation
↓
data cube.



climbing up the
concept hierarchy

or dim^n reduction.

→ Drill Down^o -

→ Reverse of roll-up

→ navigates from



less detail



more detail

→ moving down the
concept hierarchy

→ Introduce additional
 dim^n

Slice & dice^o -

Slice & dice -

Slice :-

→ performs selection
on one dimension

Dice → defines subcube
by performing selection on
two or more dimensions

Pivot (rotate) :-

→ → →

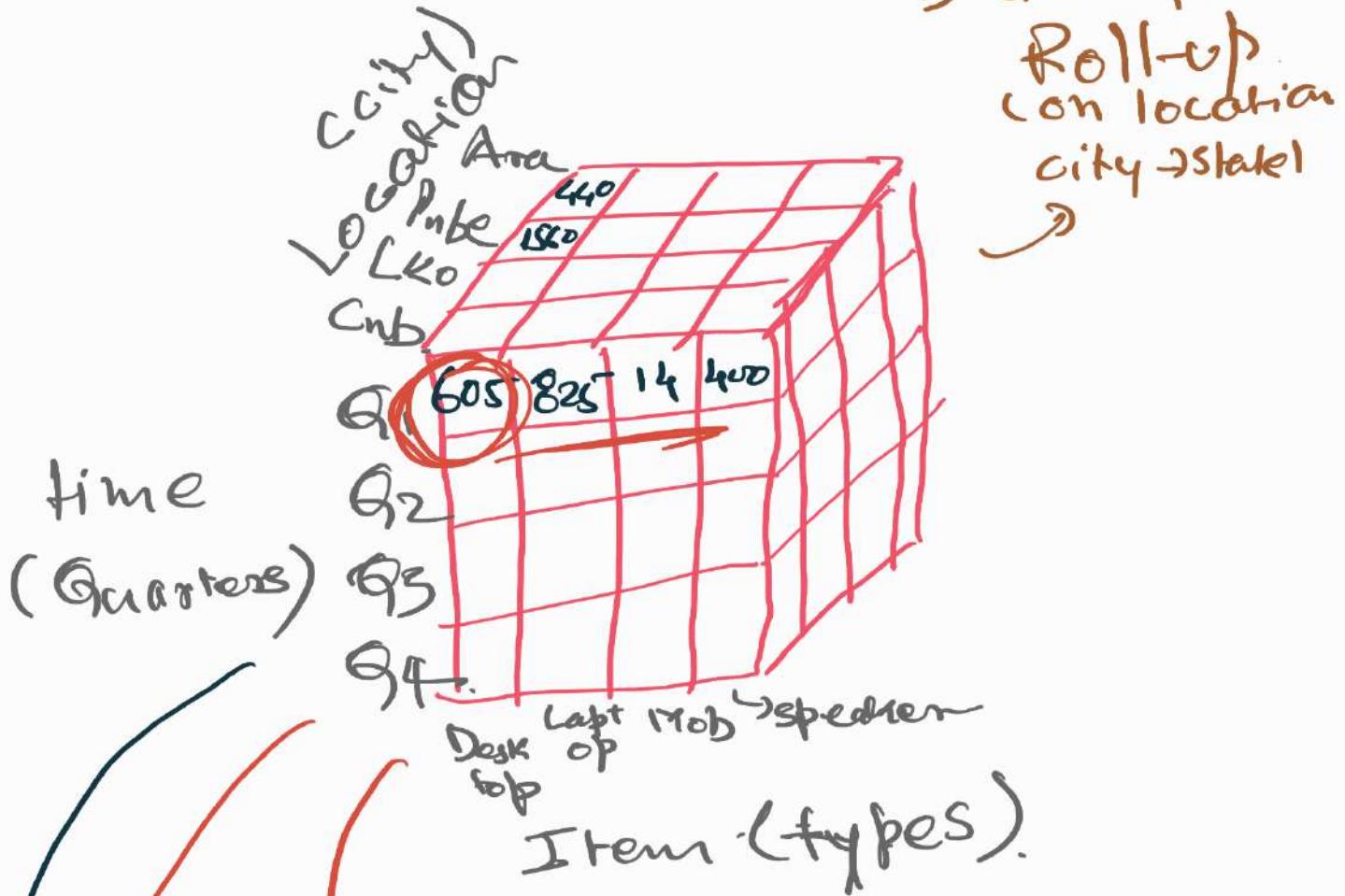
→ it rotates the
data axis in view



alternative presentation

Example :-

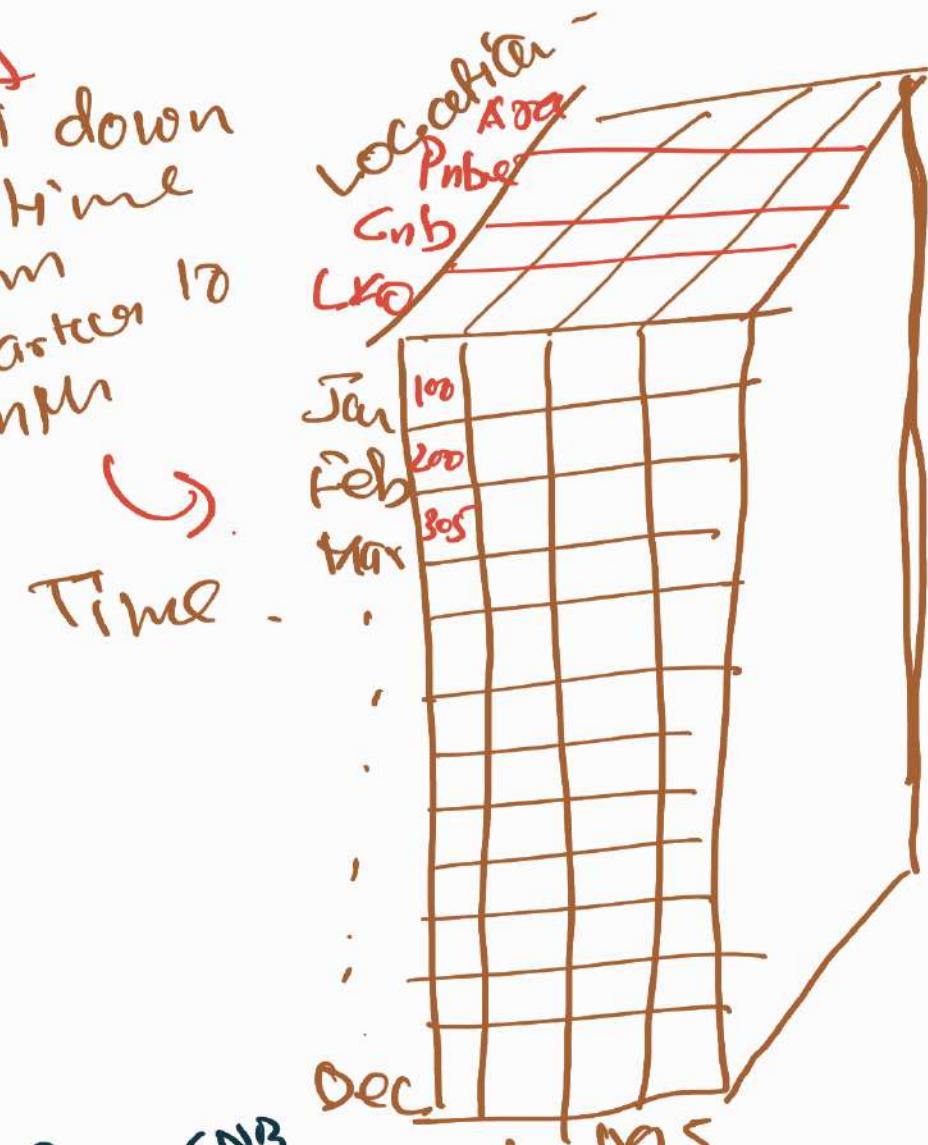
Bihar	2000	
U.P.		
Q1		
Q2		
Q3		



Drill down
on Time
(from
Quarter to
month)

Time -

Dice

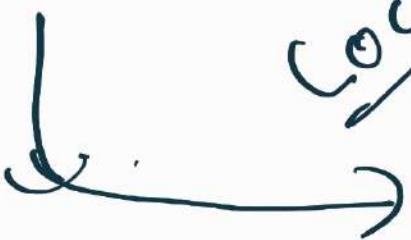


~~Location = LKO or Q1 or Q2~~ DL MS.

~~Time = Q1 or Q2~~ Item

~~Item = D or L~~ Location.

~~COCOMAN GNB~~



LKO

Q₁

Time

Q₂

D

L

Item

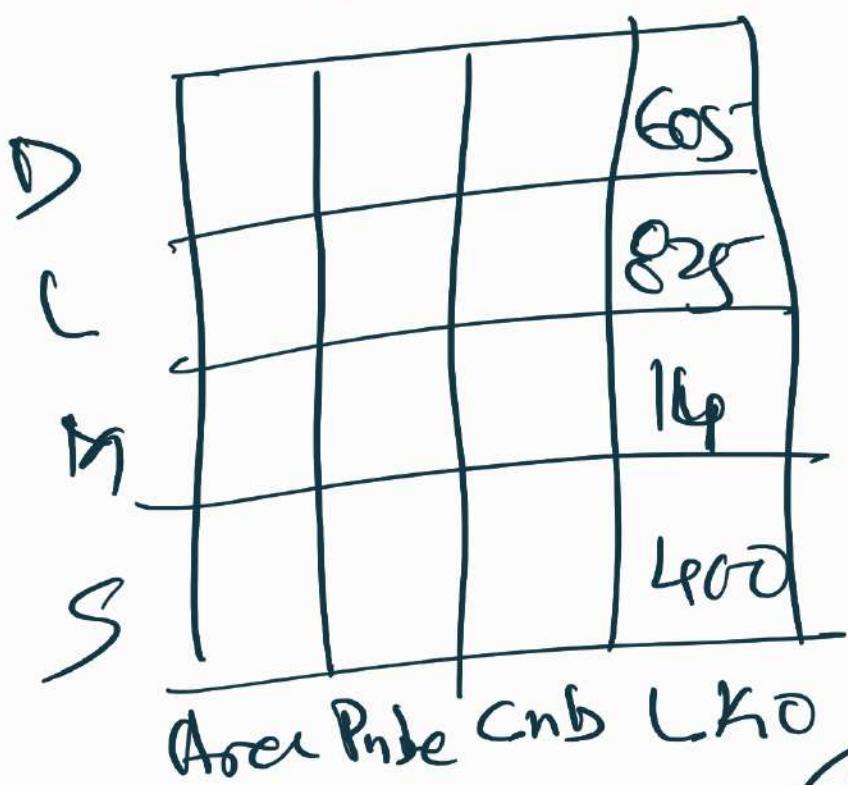
~~Slice~~

for Time = "Q1"



Area	D	L	TA S.
Pnbe	605	825	14
Card	429		
CKD	605	825	14
	D	L	TA S.

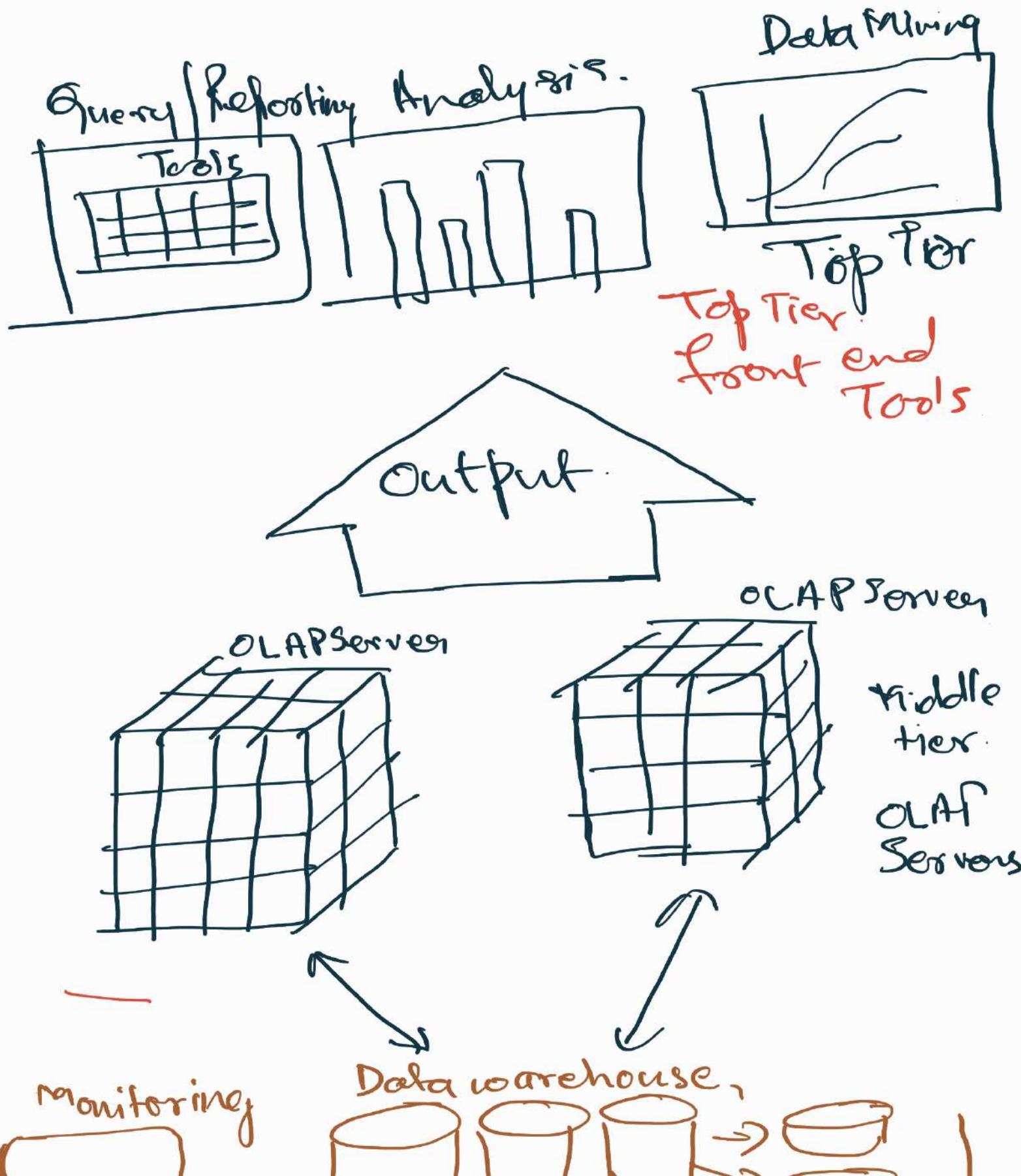
Pivot (Rotate the axes to change the view).

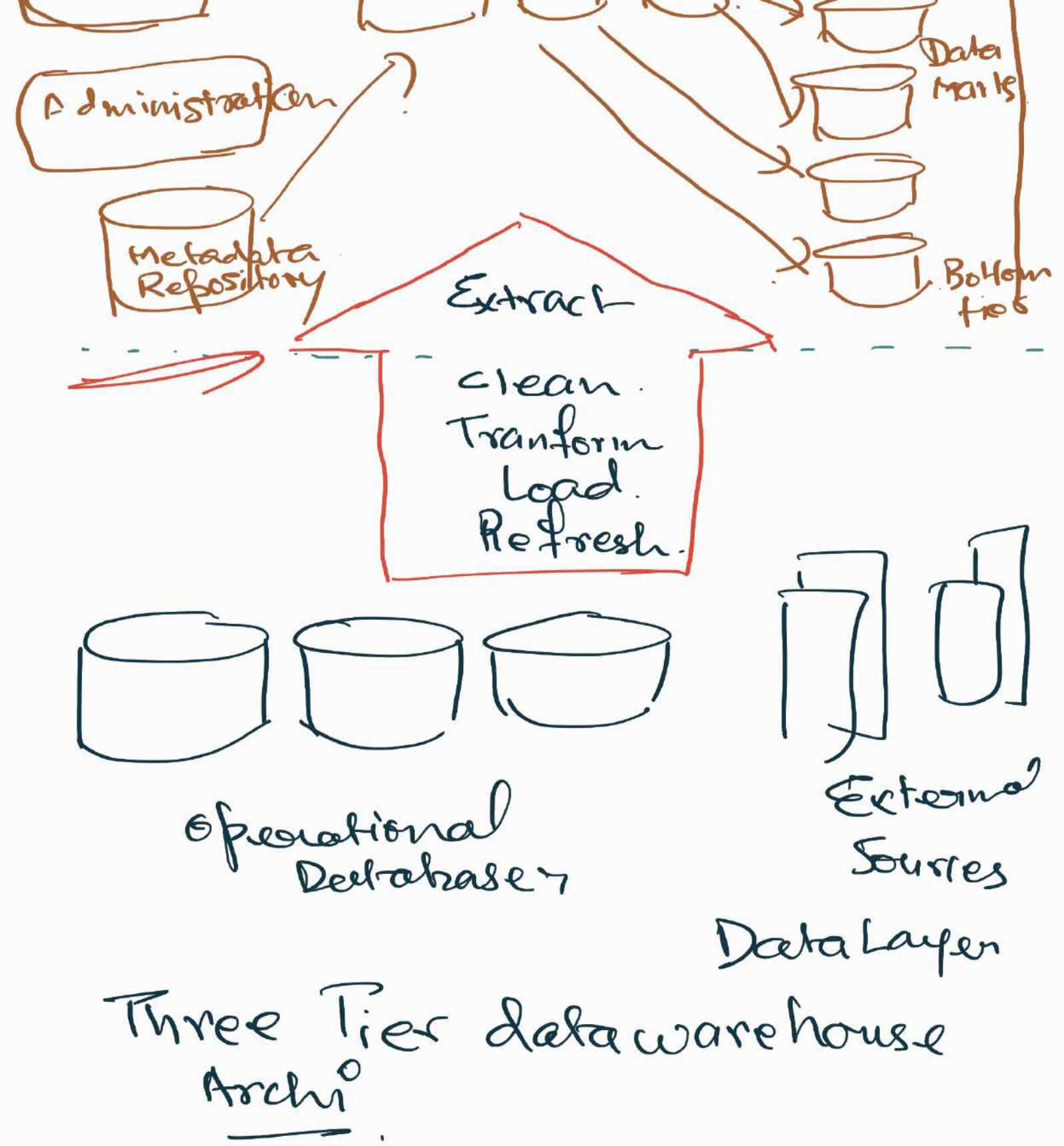


→ Han & Kamber



Three-Tier Data Warehouse Archi^o





Middle Tier -

↓
If it is an OLAP Server

~~ROLAP~~
→ MOLAP.

ROLAP (Relational OLAP)

→ extended form of relational DBMS

MOLAP → (Multidimensional OLAP)

* special purpose servers

most directly implements

"multidim" data &

operations

Top Tier

→ front end client

layer.

→ Query / reporting /
Analysis tools

Enterprise Data warehouse

- all information .
- about various subjects .
- entire organization .
- cross functional in scope
- Detailed & summarized data .
- Terabytes

Data mart —

- subset of corporate wide data .

→ specific group of users.



Independent —

→ Data directly from the external information providers / other operational system

Dependent —

→ data from the enterprise data warehouse

Frequent Pattern



→ all observations in

→ A frequent dataset frequently

Ex milk & bread.

→ Subsequence

↳ sequence that occurs frequently in the dataset-

(↳ frequent sequential pattern)

PC → Printer → Mouse Pad.

Market Basket Analysis

→ Ex. frequent itemset mining

→ Analysis customer buying habits

→ Association among diff items that ~~she~~ he places in their shopping cart

→ Retailers to develop marketing strategies.

Association Rule :-

Each item → Boolean variable.
(0/1)

Customer Basket → Boolean vector

$C_1 \rightarrow I_1, I_2, I_3, I_4, I_5, \dots, \dots, I_n$

I_1	0	1	0	1	1	0	0	0	0	0	0
-------	---	---	---	---	---	---	---	---	---	---	---

* Boolean vector can then be analyzed for

Frequently purchased
or
,, Associated items

→ Item patterns

↓
Association Rules

Ex :- Computers, Antivirus SW

Computer → Antivirus SW

[Support = 2% Confidence = 60%]

* → measuring the

Association rule

interestingness

2% support :-

2% of all the transaction under analysis show that comp & antivirus viruses are purchased together.

60% confidence :-

→ 60% of the customer who purchased of Computer also bought the

Antivirus Software -

→ minimum support threshold

→ minimum confidence threshold

A → B



